

1 Review

In the last lecture, we talked about EXP algorithm in Multi-Armed Bandit problem. In this section, we will give a review of the EXP3[1] and EXP4[1] algorithm.

1.1 EXP3 - Exponential-Weight Update algorithm for Exploration and Exploitation

EXP3 is designed for sampling action in context-free adversarial environment. The algorithm is listed as Algorithm 1 below. It is very similar to Hedge Algorithm, with two difference. First, in exp3 we only know loss for one action in a single round, leading to update only one weight. On the other hand, we are able to acquire loss for all actions in the hedge algorithm and update for all weights. The second difference is that EXP3 directly get reward from taking action, instead of sampling expert.

Algorithm 1 EXP3($\gamma \in [0, 1]$)

1: $\mathbf{w}^{(1)} \leftarrow \{\mathbf{w}_k^{(1)} = 1\}_{k=1}^K$	weights over actions
2: for $t = 1, \dots, T$ do	
3: $\mathbf{p}^{(t)} = \frac{\mathbf{w}^{(t)}}{\sum_k w_k^{(t)}}$	probability over actions
4: $k \sim \text{MULTINOMIAL}(\mathbf{p}^{(t)})$	take and draw action
5: $a^{(t)} = a_k$	
6: $\text{RECEIVE}(r^{(t)} \in [0, 1])$	get reward
7: $w_k^{(t+1)} = w_k^{(t)} e^{\gamma \frac{r^{(t)}}{p_k^{(t)}}}$	update weight
8: end for	

1.1.1 Regret Bound

Given Multi-Arm Bandit problem with K possible actions and T round, the regret bound is

$$R \leq O(\sqrt{TK \log K})$$

by setting

$$\gamma = \sqrt{\frac{\log K}{TK}}$$

Thus EXP3 is a no regret algorithm.

1.2 EXP3 Corrected

Another version of EXP3 algorithm is also presented in the lecture. The algorithm is listed as Algorithm 2. The main difference between two version is in line 3, where the corrected version generates probability through a weighted combination between weights and uniform distribution. The regret bound is exactly the same as before and the corresponding optimal γ being

$$\gamma = \sqrt{\frac{T}{2T}}$$

Algorithm 2 EXP3($\gamma \in [0, 1]$)

```

1:  $\mathbf{w}^{(1)} \leftarrow \{\mathbf{w}_k^{(1)} = 1\}_{k=1}^K$                                 weights over actions
2: for  $t = 1, \dots, T$  do
3:    $\mathbf{p}^{(t)} = (1 - \gamma) \frac{\mathbf{w}^{(t)}}{\sum_k w_k^{(t)}} + \frac{\gamma}{K}$                 probability over actions
4:    $k \sim \text{MULTINOMIAL}(\mathbf{p}^{(t)})$                                 take and draw action
5:    $a^{(t)} = a_k$ 
6:    $\text{RECEIVE}(r^{(t)} \in [0, 1])$                                 get reward
7:    $w_k^{(t+1)} = w_k^{(t)} e^{\frac{\gamma}{K} \frac{r^{(t)}}{p_k^{(t)}}}$                                 update weight
8: end for

```

1.3 EXP4 - EXP3 with Experts

EXP4 algorithm is designed for contextual Multi-Arm Bandit problem, which can be viewed as a Multi-Arm Bandit problem with N experts. The algorithm for N experts with K actions is listed as Algorithm 3. The regret bound of EXP4 is

$$R \leq \sqrt{KT \log N}$$

which suggests that the algorithm is useful when having few good experts.

Algorithm 3 EXP4($\gamma \in [0, 1], T$)

```

1:  $\mathbf{w}^{(1)} \leftarrow \mathbf{1} \in \mathbb{R}^N$                                 weights over experts
2:  $\text{RECEIVE}(\mathbf{X} \in \mathbb{R}^{N \times K})$                                 get K actions advices from N expertss
3: for  $t = 1, \dots, T$  do
4:    $\mathbf{q}^{(t)} = \frac{\mathbf{w}^{(t)}}{\sum_k w_k^{(t)}} \mathbf{X}^{(t)} \in \Delta^K$                 weighted probability of actions over N experts
5:    $k \sim \text{MULTINOMIAL}(\mathbf{q}^{(t)})$                                 take and draw action
6:    $a^{(t)} = a_k$ 
7:    $\text{RECEIVE}(r^{(t)} \in [0, 1])$                                 get reward
8:    $\hat{\mathbf{r}}^{(t)} = \frac{r^{(t)}}{q_k^{(t)}} \mathbb{I}[k = k^{(t)}] \in \mathbb{I}^K$                 reward over all arm
9:    $\mathbf{g}^{(t)} = \mathbf{X}^{(t)} \hat{\mathbf{r}}^{(t)} \in \mathbb{R}^N$                                 per expert reward
10:   $w_k^{(t+1)} = w_k^{(t)} e^{\gamma g_n^{(t)}} \forall n$                                 update weight
11: end for

```

2 Summary

2.1 Reinforcement Learning

Problem Formulation

Reinforcement Learning differs from our previous problem with its sequential feedback. The problem can be formulated as

$$\zeta, R \sim D(\zeta, R)$$

where R is the reward and ζ being

$$\zeta_i = \{(x_1, a_1), (x_2, a_2), \dots, (x_T, a_T)\} R_i \in \mathbb{R}$$

x_i is the current state and a_i is the action to take. The key difference between RL and supervised learning is that every action may affect next state; that being said, all samples in the trajectory are correlated. The goal of our learning algorithm is to learn a mapping from state x to action a

$$\pi : x \rightarrow a$$

Review of Learning Problem

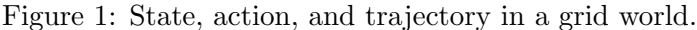
Problem	Sampled	Evaluative	Sequential
PWEA	X	X	X
OLC	V	X	X
MAB	X	V	X
C-MAB	V	V	X
RL	V	V	V

The above is the table of learning problems that we have encountered and their attributes. Comparing to all previous problems, reinforcement learning is sampled, evaluative, and sequential. The problem is sampled as for each time t it sampled the current state ζ from the distribution D . The problem is evaluative as the agent would only get reward of current action, resulting in partially observable loss. The problem is sequential as the current action will affect future state and reward.

2.2 Markov Decision Process

2.2.1 Notations

Before we talk about Markov Decision Process (MDP)[2], let's first look at an concrete example to understand terms state, action, and trajectory.



$s \in \mathcal{S}$	<i>State</i>
$a \in \mathcal{A}$	<i>Action</i>
$p(s' s, a)$	<i>State Transition Dynamic</i>
$r(s' s, a)$	<i>Reward function</i>
$p_0(s)$	<i>State prior</i>
$\pi(a s)$	<i>Policy</i>
γ	<i>Discount factor</i>

MDP is way of sequential decision making and we can decompose it to a temporal sequence of variables. Let's define a trajectory as a sequence of states and actions $\zeta = \{s_0, a_0, s_1, a_1, \dots, s_T, a_T\}$. Based on the trajectory, we consider the joint distribution of generating a trajectory $p(s_0, a_0, s_1, a_1, \dots, s_T, a_T)$ and a reward function $r(s_0, a_0, s_1, a_1, \dots, s_T, a_T)$, which is a scalar value for one trajectory. A traditional MDP factorization is as below:

$$p(s_0, a_0, \dots, s_T, a_T) = p_0(s_0) \prod_{t=0}^{T-1} p(s_{t+1}|s_t, a_t)p(a_t|s_t)$$

Return =	$r(s_0, a_0, s_1, a_1, \dots, s_T, a_T)$	or
	$r(s_0, a_0, s_1) + r(s_1, a_1, s_2) + \dots$	or
	$r(s_0, a_0) + r(s_1, a_1) + \dots$	or
	$r(s_0) + r(s_1) + \dots$	

, where $r(s)$ is the reward.

We can also look at it with a factor graphical model (Fig. 2).

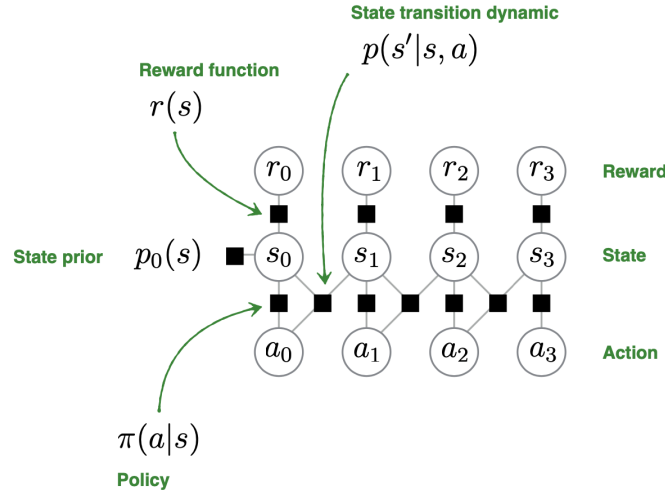


Figure 2: Graphical Model of MDP

Let's also describe the policy, state transition dynamic, and reward in the grid world (Fig. 3). Policy describes which action to take in a given state. State transition dynamic describes the probability of transition to another state. Reward function maps a state to a real value; in this example, we have the reward value finite.

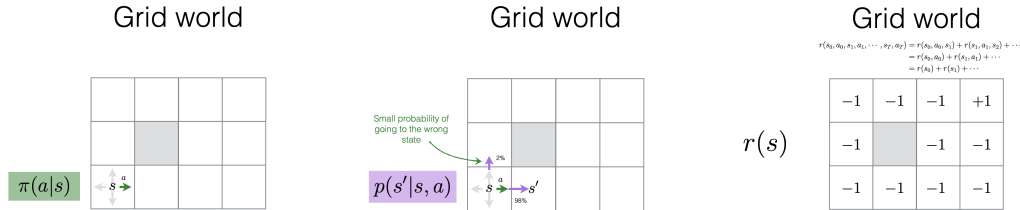


Figure 3: Policy, state transition dynamic, and reward in a grid world.

There are many inference tasks for the MDP. Fig.4 visualize the problems to solve for. Reinforcement learning, where we target at solving for policy, and inverse reinforcement learning, where we target at solving for reward function will be covered in class.

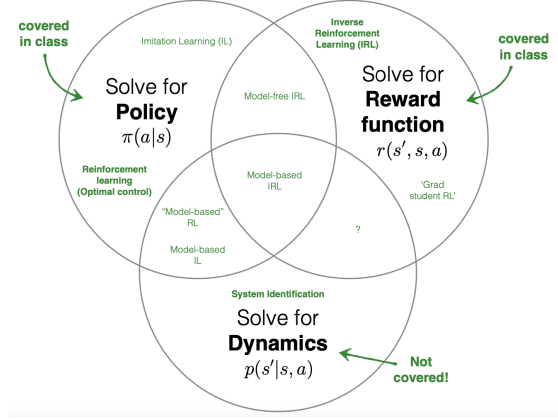


Figure 4: Inference tasks of MDP

There are many RL approaches to find the optimal policy. We can divide it into policy-based, valued-based, and hybrid. There are also many assumptions when finding the optimal policy. Two big categories are:

- Model-based
 - Dynamics known
 - Reward function known
- Model-free
 - Dynamics unknown
 - Reward function unknown
 - Learn through interaction

2.2.3 MDP Concepts

Value Function Value function is to predict the future. A type of value functions is *state value function*, which only depends on states. It is the total expected return of a trajectory in a state s .

$$V^\pi(s) = \mathbb{E}_p[r_0 + r_1 + r_2 + \dots | s_0 = s]$$

, where $r_t \triangleq r(s_{t+1} = s', a_t = a, s_t = s)$ and p is the joint distribution of the MDP. Since it is hard to deal with infinite return, we usually have *finite horizon return* and *infinite horizon discounted return*:

$$V^\pi(s) = \mathbb{E}_p[r_0 + r_1 + r_2 + \dots + r_T | s_0 = s]$$

$$V^\pi(s) = \mathbb{E}_p[\gamma^0 r_0 + \gamma^1 r_1 + \gamma^2 r_2 + \dots | s_0 = s]$$

Another type of value function is *state-action value function*. It is the total expected return of a trajectory starting in state s and taking action a .

$$Q^\pi(s, a) = \mathbb{E}_p[\gamma^0 r_0(s_0) + \gamma^1 r_1(s_1) + \gamma^2 r_2(s_2) + \dots | s_0 = s, a_0 = a]$$

The relationship between V and Q is

$$V^\pi = \sum_a \pi(a|s) Q^\pi(s, a)$$

We use the value function to be the objective function of RL. In other words, we want to find the best policy by maximizing the value function.

$$\begin{aligned} \hat{\pi} &= \arg \max_{\pi} V^\pi(s) \\ &= \arg \max_{\pi} \mathbb{E}_p \left[\sum_{t=0}^{\infty} \gamma^t r_t \mid s_0 = s \right] \end{aligned}$$

References

- [1] P. Auer, N. Cesa-Bianchi, Y. Freund, and R. E. Schapire. The nonstochastic multiarmed bandit problem. *SIAM journal on computing*, 32(1):48–77, 2002.
- [2] R. Bellman. A markovian decision process. *Journal of mathematics and mechanics*, 6(5):679–684, 1957.